

CORRECTED VERSION

(19) World Intellectual Property Organization
International Bureau



(43) International Publication Date
8 September 2000 (08.09.2000)

PCT

(10) International Publication Number
WO 00/52576 A1

(51) International Patent Classification: G06F 11/07,
11/16, 11/20, 13/14, 13/36, 12/00

(21) International Application Number: PCT/US00/03275

(22) International Filing Date: 8 February 2000 (08.02.2000)

(25) Filing Language: English

(26) Publication Language: English

(30) Priority Data: 09/261,906 3 March 1999 (03.03.1999) US

(81) Designated States (national): AE, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, CA, CH, CN, CR, CU, CZ, DE, DK, DM, EE, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM, TR, TT, TZ, UA, UG, US, UZ, VN, YU, ZA, ZW.

(84) Designated States (regional): ARIPO patent (GH, GM, KE, LS, MW, SD, SL, SZ, TZ, UG, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, GW, ML, MR, NE, SN, TD, TG).

(71) Applicant (for all designated States except US): YOTTA YOTTA, INC. (US/US); 3462 Yarrow Point Road, Bellevue, WA 98004 (US).

Published:
— with international search report

(48) Date of publication of this corrected version:
30 August 2001

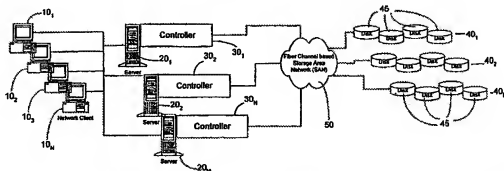
(72) Inventors; and
(75) Inventors/Applicants (for US only): BROWN, William, P. (US/US); 15431 S.E. 42nd Street, Bellevue, WA 98006 (US), MATHEWS, Michael, B. (US/US); 10725 126th Place N.E., Kirkland, WA 98033 (US).

(15) Information about Correction:
see PCT Gazette No. 35/2001 of 30 August 2001, Section II

(74) Agents: GRAY, Gerald, T. et al.; Townsend and Townsend and Crew LLP, Two Embarcadero Center, Eighth Floor, San Francisco, CA 94111 (US).

For two-letter codes and other abbreviations, refer to the "Guidance Notes on Codes and Abbreviations" appearing at the beginning of each regular issue of the PCT Gazette.

(54) Title: METHODS AND SYSTEMS FOR IMPLEMENTING SHARED DISK ARRAY MANAGEMENT FUNCTIONS



(57) Abstract: Multiple Array Management Functions (80) included in controller (30) are connected to multiple redundancy groups (40) over a storage area network (SAN), such as a fiber-channel based SAN (50). The multiple Array Management Functions share management responsibility of the redundancy groups, each of which typically includes multiple resources spread over multiple disks (45). The AMFs provide concurrent access to the redundancy groups for associated host systems. When a host requests an AMF to perform an operation on a resource, the AMF synchronizes with the other AMFs sharing control of the redundancy group that includes the resource to be operated on, so as to obtain a lock on the resource. While performing the operation, the AMF send replication data and state information associated with the resource such that if the AMF fails, any of the other AMFs are able to complete the operation and maintain data reliability and coherency.

METHODS AND SYSTEMS FOR IMPLEMENTING SHARED DISK ARRAY MANAGEMENT FUNCTIONS

5

CROSS REFERENCE TO RELATED APPLICATIONS

This application is related to U.S. Provisional Patent Application Serial No. 60/101,742, filed September 24, 1998, entitled "Methods and Systems for Implementing Shared Disk Array Management Functions," the disclosure of which is hereby incorporated by reference in its entirety.

10

BACKGROUND OF THE INVENTION

The present invention relates in general to systems and methods for eliminating bottlenecks in data storage networks, and in direct server attached storage, and more specifically to systems and methods for implementing dynamically shared redundancy group management between multiple disk array management functions.

15

The need for faster communication among computers and data storage systems requires ever faster and more efficient storage networks. In recent years, implementation of clustering techniques and storage area networks (SANs) have greatly improved storage network performance. In a typical storage network, for example, N servers are clustered together for a proportional performance gain, and a SAN (e.g., a Fiber Channel based SAN) is added between the servers and various RAID ("Redundant Array of Inexpensive Disks") storage systems/arrays. The SAN allows any server to access any storage element. However, in the typical storage network, each RAID system has an associated RAID controller that must be accessed in order to access data stored on that particular RAID system. This can lead to bottlenecks in system performance as the storage managed by a particular RAID controller can only be accessed through that RAID controller. Furthermore, if a controller fails, information maintained in the RAID system managed by the failed controller becomes inaccessible.

20

25

One solution for providing fault tolerance is to include a redundant controller in a master/slave arrangement. The master controller has primary control, and only when the master fails does the slave controller take over. This solution is very

30

inefficient, however, as the slave controller is not used until a failure in the master has occurred. Another solution is to use the master/slave controller architecture, but to split the storage array into two redundancy groups, each of which is controlled by one and only one of the two controllers (each controller is a "master" vis-à-vis the redundancy group it controls). In this manner, both controllers are operational at the same time, thereby improving the efficiency of the system. In the event one controller fails, the other controller assumes control of the failed controller's redundancy group. This solution also prevents "collisions", which occur, for example, when more than one controller tries to write data to a redundancy group. However, this solution also has some performance drawbacks. For example, the performance in such a master/slave architecture is bound by the speed of the master controller such that performance is not scalable.

Thus, it is desirable to provide techniques for implementing a peer-to-peer controller architecture solution where system performance is not bound by the speed of a given controller. Further, such a system should provide suitable fault tolerance and performance scalability.

SUMMARY OF THE INVENTION

The present invention provides such a peer-to-peer controller architecture solution for data storage management. The systems and methods of the present invention implement a novel type of RAID Array Management Function that is useful for building highly scalable disk arrays. In particular, the systems and methods of the present invention provide for sharing redundancy group management between multiple (two or more) Array Management Functions.

According to the invention, multiple Array Management Functions (AMFs) are connected to multiple redundancy groups over an interconnect medium. In one embodiment, for example, the Array Management Functions are connected to the redundancy groups over any storage area network (SAN), such as a fiber-channel based SAN. The multiple AMFs share management responsibility of the redundancy groups, each of which typically includes multiple resources spread over multiple disks. The AMFs provide concurrent access to the redundancy groups for associated host systems. When a host requests an AMF to perform an operation on a resource, the AMF synchronizes with the other AMFs sharing control of the redundancy group that includes the resource to be operated on, so as to obtain a lock on the resource. While performing

the operation, the AMF send replication data and state information associated with the resource such that if the AMF fails, any of the other AMFs are able to complete the operation and maintain data reliability and coherency.

As used herein, the terms "Array Management Function," "Redundancy Group," and "Redundancy Group Management" are defined as set forth in The RAID Advisory Board's (RAB) Handbook on System Storage Technology, 6th edition, the contents of which are herein incorporated by reference for all purposes.

"Array Management Function" (AMF) generally refers to the body that provides common control and management for one or more disk or tape arrays. An AMF presents the arrays of tapes or disks it controls to the operating environment as one or more virtual disks or tapes. An AMF typically executes in a disk controller, an intelligent host bus adapter or in a host computer. When it executes in a disk controller, an AMF is often referred to as firmware. One or more AMFs can execute in each controller, adapter or host as desired for the particular application.

"Redundancy Group" generally refers to a collection of p_extents organized by an AMF for the purpose of providing data protection. With one redundancy group, a single type of data protection is used. All the user data storage capacity in a redundancy group is protected by check data stored within the group, and no user data capacity external to a redundancy group is protected by check data within it. Redundancy groups typically include logical entities composed of many resources such as stripes, data blocks, cached data, map tables, configuration tables, state tables, etc.

"Redundancy Group Management" generally refers to the responsibilities, processes and actions of an AMF associated with a given redundancy group.

According to the invention, updates of the check data within a redundancy group are dynamically coordinated and synchronized between the various AMFs sharing the redundancy group. Such updating is facilitated using coherency and locking/unlocking techniques. Coherency and locking are typically performed as a function of a block, a group of blocks, a stripe or a group of stripes. Locking is performed dynamically using any of a variety of well known or proprietary coherency protocols such as MESI. Additionally, the coherency between the caches associated with a redundancy group and the data contained within the redundancy group is synchronized and maintained.

According to an aspect of the invention, a data storage network is provided which typically comprises a redundancy group including a plurality of resources, and two or more array management functions (AMFs) sharing access to the redundancy group. The AMFs provide concurrent access to the redundancy group for associated host
5 systems. The network also typically includes a storage area network for connecting the AMFs with the redundancy group. In a typical operation, when a first one of the AMFs desires to perform an operation on a first resource in the redundancy group, the first AMF arbitrates with the other AMFs sharing access to the redundancy group for a lock on the first resource. Thereafter, the first AMF performs the operation on the first resource and
10 concurrently sends replication data and state information associated with the first resource to the other AMFs such that if the first AMF fails while performing the operation, one of the other AMFs is able to complete the operation.

According to another aspect of the invention, a method of dynamically sharing management of a redundancy group between two or more array management
15 functions (AMFs) is provided, where the AMFs are able to concurrently access the redundancy group, which includes a plurality of resources. The method typically comprises the steps of receiving a request from a host by a first one of the AMFs to perform a first operation on a first one of the resources, synchronizing with the other AMFs so as to acquire access to the first resource, and performing the first operation on
20 the first resource.

According to yet another aspect of the invention, a data storage network system is provided, which typically comprises one or more redundancy groups, each redundancy group including multiple resources spread over multiple disks, and two or more array management functions (AMFs) sharing redundancy group management of the
25 one or more redundancy groups, wherein the AMFs are able to concurrently access the one or more redundancy groups. The system also typically comprises a storage area network for interconnecting the AMFs with the redundancy groups.

According to a further aspect of the invention, a method of reconstructing a redundancy group when one of its disks fails in a data storage network system is
30 provided. The typical network system comprising two or more array management functions (AMFs) interconnected with the redundancy group over a storage area network, wherein the AMFs all share management of the redundancy group, and wherein the AMFs are able to concurrently access the redundancy group. The redundancy group

- includes multiple resources spread over multiple disks and a replacement disk. The method typically comprises the steps of arbitrating for control of a first resource of the redundancy group by a first one of the AMFs, arbitrating for control of a second resource of the redundancy group by a second one of the AMFs, and concurrently reconstructing
- 5 the first and second resources using the replacement disk.

- According to yet a further aspect of the invention, a method of expanding a redundancy group when an extra disk is added to it in a data storage network system is provided. The network system typically comprises two or more array management functions (AMFs) interconnected with the redundancy group over a storage area network.
- 10 The redundancy group includes multiple resources spread over multiple disks. The AMFs all share management of the redundancy group, and are able to concurrently access the redundancy group. The method typically comprises the steps of arbitrating for control of a first resource by a first one of the AMFs, arbitrating for control of a second resource by a second one of the AMFs, and concurrently expanding the first and second resources
- 15 using the extra disk.

- According to an additional aspect of the invention, a method of pipelining replication of incoming host data in a data storage network system is provided. The network system typically comprises a redundancy group interconnected with two or more array management functions (AMFs) over a storage area network. The redundancy group
- 20 includes multiple resources spread over multiple disks. The AMFs all share management of the redundancy group, and are able to concurrently access the redundancy group. The method typically comprises the steps of receiving a write command by a first AMF from a host to write at least two data sets to two or more of the resources, and acquiring a lock by the first AMF on the first resource to which the first data set is to be written. The
- 25 method also typically includes the steps of writing the first data set to the first resource, and concurrently performing a first replication operation wherein replication data and state information associated with the first resource is sent to the other AMFs, such that if the first AMF fails while performing the write operation, one of the other AMFs is able to complete the write operation.

- 30 According to yet an additional aspect of the invention, a method is provided for dynamically sharing management of a redundancy group between two or more array management functions (AMFs) in a data storage system. The AMFs are able to concurrently access the redundancy group, which includes a plurality of resources.

The method typically comprises the step of determining an arbiter AMF for a first one of the resources, wherein the arbiter AMF is one of the two or more AMFs sharing management of the redundancy group. The arbiter AMF is able to grant a lock for the first resource. The method also typically comprises the steps of communicating a lock request from a first one of the AMFs to the arbiter AMF requesting a lock on the first resource, and performing an operation on the first resource by the first AMF once the lock on the first resource has been granted by the arbiter AMF.

Reference to the remaining portions of the specification, including the drawings and claims, will realize other features and advantages of the present invention.

Further features and advantages of the present invention, as well as the structure and operation of various embodiments of the present invention, are described in detail below with respect to the accompanying drawings. In the drawings, like reference numbers indicate identical or functionally similar elements.

BRIEF DESCRIPTION OF THE DRAWINGS

Figures 1 to 7 show exemplary configurations useful for providing data from one or more redundancy groups to one or more host systems using controllers sharing access to and control of redundancy groups according to the present invention;

Figure 8 shows a multiple controller configuration and the internal configuration of the controllers according to the present invention;

Figure 9 shows an operation using a general synchronization sequence according to an embodiment of the present invention;

Figure 10 shows an operation using a general replication sequence according to an embodiment of the present invention;

Figure 11a shows the flow for read operations when the redundancy group is in a normal, non-degraded mode, according to an embodiment of the present invention;

Figure 11b shows the flow for read operations when the redundancy group is in a degraded mode, according to an embodiment of the present invention;

Figure 12 shows the flow for pipelining the replication of incoming host data according to an embodiment of the present invention;

Figure 13a shows the flow for a write operation when the redundancy group is in a normal, non-degraded mode according to an embodiment of the present invention;

5 Figure 13b shows the flow for a recovery process when the AMF updating the stripe as shown in Figure 13a fails before completing the update according to an embodiment of the present invention;

Figure 14a shows the flow for a write operation when the redundancy group is in a degraded (with a failed drive) mode, according to an embodiment of the present invention;

10 Figure 14b shows the flow for a recovery process when the AMF updating the stripe as shown in Figure 14a fails before completing the update according to an embodiment of the present invention;

Figure 15 shows the flow for a background reconstruction process according to an embodiment of the present invention;

15 Figure 16 shows the general sequence flow for a background expansion process according to an embodiment of the present invention;

Figures 17a and 17b illustrate AMF communication without, and with, the message gathering techniques of the present invention, respectively;

20 Figures 18a illustrates a basic arbitration process where an AMF requests a lock for a particular resource according to the present invention;

Figure 18b illustrates the general process flow of the generalized arbitration process according to the present invention;

Figure 19 illustrates a simplified arbitration process between two AMFs in a cluster configuration for a single resource; and

25 Figure 20 illustrates exemplary resource arbitration sequences for a cluster including four AMFs according to the present invention.

DETAILED DESCRIPTION

The present invention provides for shared redundancy group management (SRGM) between multiple AMFs so that multiple AMFs can simultaneously access the same redundancy group. According to the present invention, distributed synchronization and replication techniques are used to coordinate the activities of all AMFs sharing a redundancy group and to maintain data reliability. Access to any redundancy group can be gained through any controller that includes an AMF that is sharing control of that redundancy group. The AMFs sharing a resource group are therefore peers.

- Additionally, if a redundancy group is configured for shared access and a controller fails, access to data through the failed controller is blocked, but the data on the redundancy group is still intact, protected from disk failures, and accessible from any other controller that includes an AMF which is sharing that redundancy group. Within a given controller, multiple AMFs may be present, in which case redundancy groups are preferably shared on a group by group basis. That is, some redundancy groups may be shared by a first group of AMFs, other redundancy groups may be shared by a second group of AMFs, and still other redundancy groups may not be shared at all. In preferred aspects, within a given controller there is a one-to-one association between an AMF and a redundancy group; that is an AMF manages only one redundancy group at a time. Alternately, an AMF manages multiple AMFs at the same time.

- Figure 1 shows a basic network configuration according to the present invention. As shown, a plurality of network clients 10_1 to 10_N are communicably coupled with a plurality of servers 20_1 to 20_N , each of which includes a controller 30. (In general, "N" is used herein to indicate an indefinite plurality, so that the number "N" when referred to one component does not necessarily equal the number "N" of a different component. For example, the number of clients 10 does not, but may, equal the number of servers 20 in Figure 1.) Each network client 10 is coupled to one or more of servers 20 over any of a number of connection schemes as required for the specific application and geographical location relative to servers 20, including, for example, an internet connection, any local area network (LAN) type connection, any wide area network (WAN) type connection, any proprietary network connection, etc. Each controller 30 includes one or more AMFs, and is communicably coupled with the multiple arrays 40 of disk drives 45 over an interconnect medium, such as a storage area network (SAN) 50. Preferably, SAN 50 is a fiber-channel based SAN. However, any SAN type, such as a

SCSI-based SAN, or any direct server interconnect such as a direct SCSI or FC connection may be used without departing from the spirit of the invention. Because each controller 30 has direct access to each array 40 over SAN 50, redundancy group management can be shared by all of controllers 30.

5 A fiber-channel based SAN is preferred because the fiber-channel standard is an open standard that supports several network topologies including point-to-point, switched fabric, arbitrated loop, and any combination of these topologies. Further, a number of protocols are supported by the fiber-channel standard, including SCSI, Asynchronous Transfer Mode (ATM), Transmission Control Protocol/Internet Protocol
10 (TCP/IP), High Performance Parallel Interface (HiPPI), Intelligent Peripheral Interface (IPI), and others. Moreover, Fiber-channel presently provides for data transfer speeds of up to 100MBps (200MBps duplex) at distances of up to 30 meters over copper cabling and up to 10 kilometers over fiber-optic cabling.

Figure 2 shows an example of multiple hosts, each with a controller
15 configured in a switch-based fiber-channel SAN according to the present invention. Each controller 30 is coupled to switches 55 in the SAN through two fiber-channel ports as shown. Thus, each controller 30 is in communication with all other controllers 30 and with disk array 40. Each controller 30 communicates with its host system over a PCI bus 35. Switches 55 are coupled to disk array 40 using the loop topology as shown.
20 However, it will be apparent to one of skill in the art that many loops can be supported through any of a number of switching topologies. In general, the more loops, the greater the data transfer rates that can be supported. The system redundancy as shown in Figure 2 is N-1, meaning that given N controllers (30₁ to 30_N) up to N-1 controllers can fail and no data on disk array 40 will be lost when all controllers 30 are configured to share
25 management of array 40. A controller failure for a specific host causes a loss of data availability for the specific host, but not for the entire system. Controller environmental faults, such as power supply failures are protected against in this system configuration because the data from one host system is synchronized to data on the other host systems according to the present invention as will be described in more detail below.

30 There is a recovery period associated with a controller failure. This is the time it takes for the surviving controllers to make sure that all critical data is again replicated within the cluster. If another controller fails during this recovery period, then data may be lost.

Figure 3 shows an example of multiple controllers and a single host configured in a switch-based fiber-channel SAN according to the present invention. Each controller 30 is coupled to the switches 55 in the SAN through two fiber-channel ports as shown, however, from 1 to N ports may be used as desired for the particular application.

- 5 Thus, each controller 30 is in communication with all other controllers 30 and with disk array 40 over the fiber-channel SAN. Further, each controller 30 communicates with the host system over one or more PCI buses 35. The controllers 30 are also able to communicate with each over the PCI buses 35. Switches 55 are coupled to disk array 40 using the loop topology as shown. However, it will be apparent to one of skill in the art
- 10 that many loops can be supported through any of a number of switching topologies. In general, the more loops, the greater the data transfer rates that can be supported. In this configuration, up to N-1 controllers can fail before loss of data availability to the host system. External controller 30_{N+1} is provided to protect against a failure of the host system. If the host system fails, no data will be lost on array 40 when controller 30_{N+1} is
- 15 configured to share management of array 40 with the controllers 30₁ to 30_N of the host system according to the present invention.

- Figure 4 shows an example of multiple hosts each with multiple controllers configured in a switch-based fiber-channel SAN according to the present invention. Each controller 30 is coupled to the switches 55 in the SAN through two fiber-
- 20 channel ports as shown, however, from 1 to N ports may be used as desired for the particular application. Thus, each controller 30 is in communication with all other controllers 30 and with disk array 40 over the fiber-channel SAN. Further, each controller 30 communicates with its host system over one or more PCI buses 35. The controllers 30 are also able to communicate with each over the PCI buses 35. Switches
- 25 55 are coupled to disk array 40 using the loop topology as shown. However, it will be apparent to one of skill in the art that many loops can be supported through any of a number of switching topologies. In general, the more loops, the greater the data transfer rates that can be supported. In this configuration, redundancy and synchronization exist between two or more controllers 30 within each host system. Where each host system
- 30 includes N controllers 30, up to N-1 controllers can fail before loss of data availability to the host system. If a host system fails, no data will be lost on array 40 when controllers 30 on other hosts are configured to share management of array 40 with the controllers 30 of the failed host system according to the present invention.

Figure 5 shows an example of multiple hosts each with multiple controllers configured in a fiber-channel arbitrated loop (FC-AL) SAN according to the present invention. Each controller 30 is coupled to the loop through the two fiber-channel ports as shown. Thus, each controller 30 is in communication with all other controllers 5 30 and with disk array 40 over the FC-AL. Further, each controller 30 communicates with its host system over one or more PCI buses 35. In this configuration, redundancy and synchronization exist between two or more controllers within each host system. Where each host system includes N controllers 30, up to N-1 controllers can fail before loss of data availability to the host system. Further, if a host system fails, no data will be 10 lost on array 40 when controllers 30 on other hosts are configured to share management of array 40 with the controllers 30 of the failed host system according to the present invention.

Figure 6 shows two independent redundancy groups managed by two independent controllers according to the present invention. Redundancy group A is 15 managed by controller 30_A of host system A, and redundancy group B is managed by controller 30_B of host system B. Also shown are external host system C and external host system D. According to one embodiment, the FC ports of controllers 30_A and 30_B function as both device and host channels. This allows each controller 30_A or 30_B to respond to service requests from its associated PCI bus 35 or from an external FC host 20 such as external host system C, external host system D or another controller such as controller 30_B or 30_A, respectively. In this manner, redundancy group A is made accessible to host system B and redundancy group B is made accessible to host system A. From the perspective of controller 30_A, for example, a request received from host system B to read or write data is treated as if it were received over associated PCI bus 35. 25 Similarly, external hosts systems C and D are able to access data on redundancy groups A and B by issuing read or write commands to the appropriate controller 30 over the fiber-channel SAN. Any number of host systems can be interconnected in this manner. Further, although only a single controller configuration is shown, other configurations can be used, such as the configuration as illustrated in Figure 4. The use of switches helps 30 isolate the disk arrays for performance scalability.

Figure 7 shows an example of an external RAID system including multiple controllers in passive PCI Backplane(s) configured in a switch-based fiber-channel SAN according to the present invention. The controllers 30₁ to 30_N are installed into one or

more passive PCI backplanes, and are configured to accept host commands from the FC ports and/or the associated PCI buses. Thus, external servers are able to access data on the various redundancy groups controlled by the controllers 30₁ to 30_N by issuing read or write requests to the appropriate controller 30. Although only one controller 30 is
5 required, performance and redundancy scales as more controllers are added. A FC-AL similar to the configuration shown in Figure 5 can alternatively be used if the use of switches 55 is undesired or impractical.

Figure 8 shows a multiple controller configuration and the internal configuration of controllers 30 according to the present invention. One or more of the
10 controllers 30₁ to 30_N shown in Figure 8 may be located in a separate host or on passive PCI backplanes. For example, each controller 30 may be located in a separate host system, or each of multiple host systems may include one or more of the controllers 30. PCI host connection 60 provides a connection path for receiving and processing commands from host systems and for providing inter-controller link (ICL) services with
15 other controllers. Fiber-channel (FC) host connection 65 provides a connection means for processing and receiving commands from host systems and for providing ICL services with other controllers. In preferred aspects, each controller includes two physical FC ports (not shown in Figure 8, but see Figures 2 through 7 for examples), both of which are used for disk drive access, receiving and processing host commands and ICL services. It
20 will be apparent to one skilled in the art that each controller can include from 1 to N FC ports as desired for the particular application.

Each controller 30 includes one or more virtual disk ports 70 each of which provides access to a virtual disk 75. Virtual disks 75 are basically partitions of an array. (A "Redundancy Group" is generally synonymous with "Array"). Each array may
25 be partitioned into as many virtual disks 75 as desired. Each virtual disk is associated and controlled by one or more associated AMFs 80. Many virtual disk ports 70 can exist for the same virtual disk 75, but each must exist on a separate controller. For example, as shown in Figure 8, virtual disk ports 70YR₁ and 70YR_N associated with virtual disk YR are provided to hosts on controller 30₁ and controller 30_N, respectively. Each virtual disk
30 port YR provides access to the same virtual disk YR. Virtual disk YR is a partition of array Y, the control and management of which is shared by AMFs 80Y₁ and 80Y_N.

Virtual disk ports can be added to a controller at any time by creating a virtual disk and assigning an IO port address to the virtual disk port. Generally, a virtual

disk must exist before a virtual disk port is created and assigned to it, but the creation of a virtual disk is not coincident with the creation of a virtual disk port. Typically a virtual disk port is created right after the redundancy group is created. Virtual disk ports can then be created at any time thereafter, but the creation of the virtual disk is only done once. Virtual disk ports can also be deleted at any time. All host operations in progress on the port are allowed to complete. While these operations are completing, new host operations are rejected, for example, by returning a not_ready status signal to the host.

Given that multiple AMFs 80 are present on each controller 30, redundancy groups are shared on a group-by-group basis. As shown in Figure 8, there are three groups of AMFs, Y, G and T, each sharing an array. Associated AMFs on different controllers synchronize and replicate with each other as will be described in more detail below. For example, as shown in Figure 8, AMF 80Y₁ synchronizes and replicates with AMF 80Y_N (and with any other AMF associated with array Y, e.g., AMF 80Y₂ (not shown)). Similarly, AMF 80G₁ synchronizes and replicates with AMF 80G_N, and AMF 80T₁ synchronizes and replicates with AMF 80T_N. Additionally, virtual disk ports on one controller synchronize and replicate with related virtual disk ports on other controllers.

Distributed Synchronization and Replication

Synchronization and replication ensures that the operations performed by the different AMFs sharing a redundancy group (array) do not destructively interfere with each other (e.g., "collisions" do not occur). Synchronization requires that any AMF which needs to access a resource associated with a shared redundancy group arbitrate with the other AMFs for access rights (lock) on the resource before using it. Arbitration is accomplished by sending arbitration messages between the AMFs over the PCI and/or FC ICL links.

Figure 9 shows a general synchronization sequence for an operation according to an embodiment of the present invention. In step 110, the operation is started. For example, a host may send a request that requires data be written to a particular resource. At step 120, the AMF determines whether it already has a lock on the desired resource. If not, the AMF arbitrates with other AMFs for a lock on the desired resource in step 130. Once a lock has been acquired (or it is determined that the AMF already has the lock), the desired operation is performed on the resource by the AMF in step 140. Once a lock is acquired by an AMF, it is preferably not released until another

AMF needs the lock (i.e., another AMF arbitrates for the lock) to help cut shared redundancy group management (SRGM) overhead in many applications. According to one embodiment, a first-come-first-served type arbitration scheme is used, but a priority based, or any other arbitration scheme can be used. As will be described in more detail below with reference to Figures 18, 19 and 20, arbitration typically involves making a request to use a resource to a resource controller (typically software, but sometimes hardware based). The resource controller grants access to the resource based on the arbitration algorithm used. If a resource is in use when the request is made, the requester waits until the resource is free and all other requesters ahead of it have been serviced.

Each AMF is able to execute many types of operations on a redundancy group, including, for example, host reads, host writes, background writes, regeneration's, reconstruction's, online expansion, parity scrubbing, etc. An extended sequence of such operations is termed a "process". Examples of processes include reconstructions, online expansion, and parity scrubbing. All AMF operation types require synchronization arbitration in order to proceed. Once an operation on a resource is completed by a given AMF, other AMFs are able to use the resource.

Synchronization is preferably performed at the operation level as opposed to the process level. That is, for each operation to be executed, the basic synchronization sequence shown in Figure 9 is performed. For a process wherein some function must be performed on the entire redundancy group (such as a reconstruction), the processing is broken up into a sequence of operations. Where each operation operates on a different resource, arbitration for synchronization of those resources required for one operation is done independently of the resources required for other operations in the process. Using synchronization at the operation level instead of the process level allows AMFs to share resources while a process is in progress. If synchronization were performed at the process level instead of the operation level, some AMFs would have to wait until the entire process is complete before they could use the resources, thereby resulting in host timeouts.

Replication accommodates AMF failures. Resources and their state information are replicated so that if an AMF fails the resource and its state information is available via another AMF that has a copy of the resource and its state information. During the update (writing) of a certain resource, a copy of the modified resource and/or the resource's operation state is sent to other AMFs sharing the resource. These other

AMFs are called replication partners. For example, referring to Figure 8, AMF 80Y₁ and AMF 80Y_N are replication partners as each share control of Array Y. The replicated information is used by the replication partners to complete the operation in the event that the AMF updating the resource fails during the operation.

5 Figure 10 shows a general replication sequence for an operation according to an embodiment of the present invention. The start of the sequence is the basic synchronization sequence as shown in Figure 9. In step 210, the operation is started. For example, a host may send a request that requires writing data to a particular resource. In step 220, the AMF determines whether it already has a lock on the desired resource. If
10 not, the AMF arbitrates with other AMFs for a lock on the desired resource in step 320. Once a lock has been acquired the operation can be performed. As shown, the operation performed by the AMF is broken into a number, i, of steps. In step 240₁, the replication data and state information associated with the resource and the first operation step is sent to each replication partner. In step 250₁, the first step of the operation is performed. As
15 shown in Figure 10, subsequent operation steps 250₂ to 250_i are performed in sequence, as are the replication steps 240₂ to 240_i. In each of replication steps 240₁ to 240_i, the replication information is sent to the replication partners associated with the given resource. In general, N-1 concurrent AMF failures are accommodated if N copies of a resource and its state information exist within the AMF cluster (i.e., those AMFs sharing
20 the resource), where N is defined as the replication dimension. Thus, replication information is sent to the N-1 replication partners associated with the given resource.

Replication can be targeted to specific replication groups or specific AMFs. In one embodiment, N-way replication is performed without defined replication groups. In this embodiment, replication takes place with any N-1 AMFs in the cluster
25 that are sharing the resource being replicated. Alternately, replication is performed with N-1 other replication groups. A replication group is a group of AMFs that replicate critical data to AMFs in replication groups other than their own. An example of this is a set of controllers, each controller including one or more AMFs, at one physical site and another set of controllers at another physical site. Another example is a set of controllers
30 inside a host system and another set external to a host. Using replication groups helps ensure that if one group of controllers all fail, the other group(s) have the information necessary to maintain data reliability.

The replication can be targeted to specific replication groups or specific AMFs. A given replication group preferably replicates with any AMF outside of the replicating AMFs replication group. However, the set of replication groups to which a given replication group replicates may be specified by an operator. Also, synchronous replication is the preferred replication mode. In the synchronous replication mode, completion status information for an operation is returned to a host after all replication targets have received the replication data. Alternate replication modes include asynchronous replication and pseudo-synchronous replication. In the asynchronous replication mode, completion status information for an operation is returned to a host before replication data has been transmitted. In the pseudo-synchronous replication mode, completion status information for an operation is returned to a host after the replication data has been transmitted, but before all replication targets have acknowledged receiving the data.

According to the present invention, multiple AMFs are able to read a resource concurrently. That is, multiple read locks can be outstanding at any time. However, only one AMF at a time writing to the resource is accommodated. Allowing multiple AMFs to read a resource at the same time significantly improves read performance. If an AMF is writing to the resource, no other AMF is permitted to read the resource. The synchronization protocol allows and enforces this aspect.

Figure 11a shows the general sequence flow for a read operation when the redundancy group (RG) is in a normal, non-degraded mode according to an embodiment of the present invention. "Non-degraded" generally refers to the case where all drives in the redundancy group are operational, whereas "degraded" generally refers to the case where one or more drives in the redundancy group have failed. In step 310, the read operation is started. For example, the AMF receives a request from a host to read a particular resource. In step 320, a lock on the particular resource is required. This is basically the same as steps 120 and 130 of Figure 9. In preferred aspects, multiple locks can be outstanding. This enables multiple AMFs to read a resource concurrently.

Figure 11b shows the general sequence flow for a read operation when the redundancy group (RG) is in a degraded mode according to an embodiment of the present invention. In step 310, the read operation is started. For example, the AMF receives a request from a host to read a particular resource. In step 320, a lock on the particular resource is required. Once the lock is acquired, in the degraded read mode, the AMF

reads the data and parity from the particular resource at step 340, and regenerates any missing data at step 350. In step 360, the data (regenerated) is transferred to the host that issued the read request.

Figure 12 shows the general sequence flow for replicating incoming host data in a pipelined manner for a write operation according to an embodiment of the present invention. Pipelining of the replication data helps to minimize replication latency. In step 410, the operation is started. For example, a host issues a write command to write one or more blocks of data to one or more resources. In step 420, the host command is received from the host and parsed. The host command is processed as a sequence of data block sets. In step 430, the appropriate lock for the starting set of blocks is acquired. In step 440, the starting block set is transferred to the AMF from the host. In step 450, the block set replication is started for the starting set of blocks. The AMF does not wait for the block set replication to complete; the AMF immediately determines whether any more sets of blocks need to be processed at step 460. If so, the AMF immediately starts acquiring the appropriate lock to get the next set of blocks in step 430, and repeats steps 440, 450 and 460 for the next block set. If all block sets have been received and processed, the AMF waits for all replication operations to complete in step 470. When each operation is complete the AMF sends status to the host in step 480.

Figure 13a shows the general sequence flow for a write operation when the redundancy group (RG) is in a normal, non-degraded mode according to an embodiment of the present invention. In step 510, the operation is started. For example, a host issues a write command to write data to a resource. Step 520 is the process of acquiring the synchronization lock for the resource required as shown in Figure 9. In this example the resource is a stripe write lock, but it may also be any other lock as required by the particular operation. After the stripe write lock is acquired, the AMF reads the old data and parity from the RG's disks in step 530.

Concurrent with the disk read operation of step 530, the AMF sends a state notification signal to its replication partners for this resource, in step 540. The replication partners include all other AMFs to which this AMF replicates state information and data for the particular resource. The number of replication partners is equal to $N-1$ where N is the replication dimension. Preferably, the replication dimension N is from 1 to 8, but N may be any number as desired. The state notification signal is a 'begin update' type signal, which tells the replication partners that a stripe update has started. The replication

partners need to know this information because they will be responsible for cleaning up in the event the writing AMF fails before completing the operation.

Once the 'begin update' notification step 530 and the read of old data and parity step 540 are complete, the AMF writes the new data to the RG member disks in step 550. Concurrent with the new data write step 550 is the generation of the new parity in step 560. Thereafter, in step 570 the AMF writes the new parity to the RG member disks. Once the parity write operation is complete, the AMF sends an 'end update' notification to its replication partners in step 580. Upon receiving this notification, the replication partners release their state data associated with the stripe update.

Figure 13b shows the general sequence flow for a recovery process when the AMF updating the stripe as shown in Figure 13a fails before completing the update according to an embodiment of the present invention. When the AMF updating the stripe fails before completing the update, the replication partners assume the responsibility of recovering from the failed update operation. In step 610, the recovery operation begins when one or more of the replication partner AMFs either detects a failure or is notified of a failure. The replication partners arbitrate for ownership of the stripe lock in step 620. The AMF that wins the arbitration (the recovery AMF) is responsible for executing recovery of the failed update operation.

Failure notification typically comes from the inter-controller link (ICL) component. If a controller fails, the AMFs on that controller lose communication with the other AMFs they were sharing the redundancy group with. The ICL periodically sends a 'ping' message to all the other AMFs it is sharing the redundancy group with. If any of these AMFs fails to respond to the ping message, then the AMF that sent the ping message assumes the AMF has failed and begins recovery action. Recovery is also triggered if the ICL encounters a transport failure when sending synchronization or replication messages to the destination AMF.

The recovery process includes two basic steps: recalculation of the stripe parity and rewriting the data. In step 630, the recovery AMF reads all the data for the stripe segments affected by the failed update operation. Concurrent with the data read step 630, the recovery AMF assigns one or more new replication partners and sends a 'begin update' notification to all of its replication partners in step 640. When the data read operation is complete, the recovery AMF generates new parity in step 650. This new

parity calculation does not include the new data. It is simply a regeneration of parity for data on the RG member disks.

Once the parity generation step 650 and 'begin update' notification step 640 are complete, the recovery AMF writes the new parity to RG member disks in step 660. After the parity write step 660 is complete, the recovery AMF sends an 'end update' notification to the replication partners in step 670. After some time interval, the cache write back scheduling algorithm causes one of the replication partners to write the new data to the RG member disks in step 680, which is a normal (non-recovery mode) stripe update operation as shown in Figure 13a. The caching functionality is a part of the AMF.

Figure 14a shows the flow for a write operation when the redundancy group (RG) is in a degraded (with a failed drive) mode, according to an embodiment of the present invention. This sequence is similar to that of the non-degraded case shown in Figure 13a, with the inclusion of regeneration step 744 and replication step 746 as will be described below.

In step 710, the operation is started. For example, a host issues a write command to write data to a resource. Step 720 is the process of acquiring the synchronization lock for the resource required as shown in Figure 9. In this example the resource is a stripe write lock, but it may also be any other lock as required by the particular operation. After the stripe write lock is acquired, the AMF reads the old data and parity from the RG's disks in step 730.

Concurrent with the disk read operation of step 730, the AMF sends a state notification signal to its replication partners for this resource, in step 740. The replication partners include all other AMFs to which this AMF replicates state information and data for the particular resource. The state notification signal is a 'begin update' type signal, which tells the replication partners that a stripe update has started. The replication partners need to know this information because they will be responsible for cleaning up in the event the writing AMF fails before completing the operation. Once the read of old data and parity step 540 is complete, the AMF regenerates the data that was on the failed disk in step 744. In step 746, the old data, including regenerated data, is replicated to the replication partners. Replication of this data to the replication partners is necessary for recovery in the event the updating AMF fails before completing the operation.

Once replication step 746 is complete, the new data is written to the RG member disks in step 750. Concurrent with the new data write step 750 is the generation

of the new parity in step 760. Thereafter, in step 770 the AMF writes the new parity to the RG member disks. Once the parity write operation is complete, the AMF sends an 'end update' notification to its replication partners in step 780. Upon receiving this notification, the replication partners release their state data associated with the stripe update.

Figure 14b shows the general sequence flow for a recovery process when the AMF updating the stripe as shown in Figure 14a fails before completing the update according to an embodiment of the present invention. This scenario differs from the non-degraded recovery case shown in Figure 13b in that the recovery AMF uses the replicated old data to reconstruct the RG disk segments that were affected by the updating AMF's failure. When the AMF updating the stripe fails before completing the update, the replication partners assume the responsibility of recovering from the failed update operation. In step 810, the recovery operation begins when one or more of the replication partner AMFs either detects a failure or is notified of a failure, for example by a host. The replication partners arbitrate for ownership of the stripe lock in step 820. The AMF that wins the arbitration (the recovery AMF) is responsible for executing recovery of the failed update operation.

In step 830, new parity is generated from the old data supplied by replication step 746 of Figure 14a. Concurrent with step 830, the recovery AMF assigns one or more new replication partners and sends a 'begin update' notification to all of its replication partners in step 840. In step 850, the old data is written to the disks. In step 860, the replication partners are informed that the old data has been written back to the disks. The replication partners can now discard their copies of the old data. After this event, the recovery sequence is the same as for the non-degraded recovery sequence. Specifically, the new parity is written to the RG member disks in step 870. After the parity write step 870 is complete, the recovery AMF sends an 'end update' notification to the replication partners in step 880. After some time interval, the cache write back scheduling algorithm causes one of the replication partners to write the new data to the RG member disks in step 890, which is a normal (non-recovery mode) stripe update operation as shown in Figure 13a.

Distributed Reconstruction

Figure 15 shows the general sequence flow for a background reconstruction process, according to an embodiment of the present invention. Each operation is started in step 910, and the appropriate lock is acquired in step 920, in this case a stripe lock. In step 930, the AMF reads the data and parity for the stripe. In step 940, the AMF regenerates missing data, and in step 950 writes the data to the replacement disk. After writing the data, the AMF updates the map table to reflect the fact that blocks that originally mapped to the failed disk now map to the replacement disk in step 960. The map table maps host logical blocks to RG member disks and blocks on the disks. In step 970, it is determined whether there are any more stripes to reconstruct. This sequence repeats as long as there are more stripes that need to be reconstructed onto the replacement disk.

One advantage of the shared redundancy group management (SRGM) techniques of the present invention is that SRGM allows for distributing reconstruction processing to all AMFs sharing a RG. This results in faster reconstruction times and a general reduction in the increased response times normally encountered by a host during reconstruction. Distributed reconstruction is accomplished by having a subset (or all) of the AMFs sharing a RG coordinate for which stripes each will be responsible for reconstructing. If any of these AMFs fail or shutdown during reconstruction, the remaining AMFs re-arbitrate reconstruction responsibilities. For example, suppose there are N AMFs sharing a redundancy group that needs reconstruction. These AMFs talk to each other (by sending messages) and determine which ones are to participate in the reconstruction, e.g., a subset of N, denote by M. These M AMFs determine reconstruction responsibilities by determining which AMFs will reconstruct which stripe(s). This can be determined by any algorithm. In one embodiment of the invention, for example, each AMF_i reconstructs stripe(s) $\text{Mod}(S/i)$, where S is the stripe number.

Distributed RG Expansion

Just as reconstruction is preferably distributed among a subset of the AMFs sharing a RG, so is RG expansion. RG expansion is the addition of drive members to an existing RG. A unique advantage of SRGM is that it allows expansion processing to be distributed to all AMFs sharing a RG. This results in faster expansion times and a

reduction in the increased response times normally encountered by a host during expansion.

Distributed expansion is accomplished by having a subset (or all) of the AMFs sharing a RG arbitrate for which stripes they will be responsible for expanding.

- 5 If any of these AMFs fail or shutdown during expansion then the remaining AMFs re-arbitrate expansion responsibilities. For example, suppose there are N AMFs sharing a redundancy group that needs expansion. These AMFs talk to each other (by sending messages) and determine which ones are to participate in the expansion, e.g., a subset of N, denote by M. These M AMFs determine expansion responsibilities by determining
- 10 which AMFs will expand which stripe(s). This can be determined by any algorithm. In one embodiment of the invention, for example, each AMF_i expands $stripe(s) \text{ Mod}(S/i)$, where S is the stripe number.

- Figure 16 shows the general sequence flow for a background expansion process according to an embodiment of the present invention. The process is started in
- 15 step 1010, and the appropriate stripe lock is acquired in step 1020. The expansion case is different from the preceding examples in that multiple locks must be acquired. An expansion operation will involve 2 or more stripes. One stripe is the expansion stripe whose stripe width will be increased from W to W+N. The other stripes involved are stripes containing host data that will be migrated from those stripes to the expanded
- 20 stripe.

- In step 1030, the data on the stripe to be expanded is read. In step 1040, the data is replicated so that if the operation fails before completion the replication partners will be able to clean up after the failure and continue the expansion process. Concurrent with steps 1030 and 1040, the source data stripe containing data to be
- 25 migrated to the expanded stripe is read in step 1045. After steps 1040 and 1045 are complete, the AMF notifies its replication partners that it is beginning the expansion stripe update in step 1050. Concurrently, in step 1055, the AMF generates parity information for the expanded stripe.

- After the notify begin update has completed, the data for the expanded
- 30 stripe is written to the disks in step 1060. Once the parity generation step 1055 and the notify begin update step 1050 are complete, the parity is written to the disks in step 1070. Once steps 1060 and 1070 are complete, the AMF notifies its replication partners that the update is complete in step 1080. The replication partners then update their map tables to

reflect the increased stripe width and migrated host data. They also discard the data replicated in step 1040. The map table maps host logical blocks to RG member disks and blocks on the disks.

- 5 In step 1090, it is determined whether any more stripes are to be expanded by the AMF. If so, the sequence is repeated. This sequence repeats as long as there are more stripes that need to be expanded to utilize the capacity of the new RG member disks. Note that this is process - what makes it a process is the looping that results in step 1090. Steps 1020 through 1090 comprise an operation.

10 Message gathering

- One issue of concern with SRGM is the processing overhead and IO channel loading associated with the synchronization and replication used to support SRGM. To facilitate the synchronization and replication techniques, ICL (Inter-Controller Link) message gathering is preferably used. Message gathering is, generally,
15 an algorithm that combines many small messages destined for a particular cluster node (i.e. a controller, in which may reside many AMFs) into one big message packet and sends it as one message to the particular node. This dramatically reduces processing overhead and IO channel loading, and contrasts with the approach of sending individual messages to a cluster node.

- 20 Figure 17a illustrates AMF communication without the message gathering techniques of the present invention. As shown, a collection of AMFs 1100 and an Inter-Controller Link (ICL) entity 1105 compose a SRGM node 1110. A node is typically a hardware entity such as a controller. ICL 1105 is a software entity that is responsible for routing synchronization and replication messages 1120 from one AMF to another. As
25 shown in Figure 17a, only one of many similar nodes is shown as being connected to the SAN 1130.

- The AMFs 1100 within node 1110 are sending and receiving synchronization and replication messages with other AMFs on other nodes that share the same redundancy group. Each AMF within node 1110 generates independent streams of
30 synchronization and replication messages, all destined for one or more other nodes on SAN 1130. The messages being sent or received by a particular AMF are independent of the messages being sent or received by other AMFs on the same node. As shown in Figure 17a, three AMFs 1100 are sending a total of nine messages 1140 to AMFs on

other nodes. Without message gathering, ICL 1105 has to send nine messages to other nodes. Also, without message gathering, all synchronization and replication messages generated by all AMFs within a SAN node are processed and sent through the SAN individually. Each message takes a fixed amount of processing overhead, regardless of size.

Figure 17b illustrates AMF communication with the message gathering techniques of the present invention. Message gathering is where many smaller messages destined for a particular node are packed together to form one larger message. This larger message can be sent over SAN 1130 as one message and then unpacked on the receiving node back into the individual messages. For example as shown, the nine messages 1120 are destined for three different nodes. In this example, then, if message gathering is used, ICL 1105 only needs to send three messages 1150 – one for each node (not counting itself). ICL 1105 takes on the responsibility of packing and unpacking individual AMF messages.

Arbitration and Distributed Resource Arbitration

Figures 18a illustrates a basic arbitration process where an AMF requests a lock for a particular resource according to the present invention. AMF 1200 and AMF 1210 each request a lock on a particular resource, and the lock requests are queued in an arbitration queue 1205. The arbitration process for an AMF begins when a request is placed in arbitration queue 1205. The requests are processed in some order such that all requests are satisfied in priority order. The request queue priority is established through any well known algorithm (e.g. FIFO, LIFO). Each requesting AMF must wait until its request is processed to obtain the lock. Each AMF obtains a lock on the resource at successful completion of the arbitration process. An AMF fails to lock the resource if arbitration fails.

Figure 18b illustrates the general process flow of the generalized arbitration process according to the present invention. Arbitration involves coordinating the resource lock transfer between two AMFs: the requesting AMF 1225 and the AMF 1230 with the resource lock. AMF 1225 sends a Request Lock message to an arbitrator 1220 (the entity performing the arbitration process), which queues the message until a time defined by its priority algorithm. When appropriate, arbitrator 1220 processes the request by issuing a Release Lock message to AMF 1230 which currently has the

resource lock. Once it has reached a point where it can safely discontinue use of the resource, AMF 1230 releases the lock and notifies arbitrator 1220 that the lock is released. Arbitrator 1220 then signals requesting AMF 1225 that it has been granted the resource lock. AMF 1225 continues to hold the lock until arbitrator 1220 calls for it to release the resource.

Optimizations of the arbitration queue are possible when one or more AMFs request a read lock for a particular resource. In preferred aspects, the arbitration process simultaneously grants read locks in these situations, as long as command ordering is preserved.

An AMF (or controller) manages the arbitration process for a resource within the redundancy group cluster. This AMF is known as the resource arbitrator. Assignment of the arbitrator for a specific resource can be accomplished using any of multiple methods (e.g. single arbitrator for all resources, load balancing assignment, etc.). The preferred methods for arbitration assignment according to the present invention are based on the number of controllers and the resource range. For cluster configurations with one or two AMFs, the assigned arbitrator is the last AMF with a Write Lock. For cluster configurations with greater than two AMFs, arbitrator assignment is based upon a fixed mapping algorithm based upon the number of AMFs and the resource location (e.g. stripe ID, block ID, etc.).

Figure 19 illustrates a simplified arbitration process between two AMFs in a cluster configuration for a single resource. First AMF 1300 (AMF #1) issues a Read Lock request 1320 for a resource to second AMF 1310 (AMF #2), which currently has a Write Lock on the resource. AMF#2 issues a Grant Lock (read) message 1330 to AMF #1 indicating that a resource lock has been granted. AMF #1 now has read access to the resource. The sequence continues when AMF #1 issues a Write Lock request 1340 to AMF #2. AMF #2 responds with a Grant Lock (write) message 1350. Continuing, AMF #1 issues a Read Lock request 1360, and since AMF #1 already has a Write Lock, it handles its own arbitration and demotes the Write Lock to a Read Lock. Next, AMF #2 has no locks on the resource at this time, so it does not need to be notified. AMF #2 issues a Read Lock request 1370 to AMF #1, which responds immediately with a Grant Lock (read) message 1380 since the resource supports multiple read locks. For the two-AMF cluster configuration, arbitration with the other AMF is not required when the AMF

with a Write Lock is requesting a Read Lock. In all other cases, the second AMF must arbitrate the lock.

Figure 20 illustrates exemplary resource arbitration sequences for a cluster including four AMFs according to the present invention. The preferred arbitrator assignment method for clusters containing three or more AMFs is to select the arbitrator using a fixed mapping algorithm. This has the effect of permanently associating an arbitrator with a single AMF. In the event of AMF resource arbitration failure, the resource arbitrator is reassigned according to the mapping algorithm.

Referring to the Figure 20, first AMF 1400 (AMF #1) issues a write lock request 1420 to the resource X arbitrator on third AMF 1410 (AMF #3). The arbitrator on AMF #3 issues a release lock request 1422 to second AMF 1405 (AMF #2) to release its lock on the resource X. While the lock request 1420 for resource X is processing, AMF #1 issues a read lock request 1424 for resource Y. Fourth AMF 1415 (AMF #4) is the assigned arbitrator for resource Y. AMF #4 immediately grants a read lock 1426 since no other AMFs currently have a lock. Next, AMF #4 issues a write lock request 1428 for resource X, which is queued by the arbitrator on AMF #3 since it is currently processing write lock request 1420.

AMF #2 sends a lock released message 1430 to AMF #3, which sends a grant lock (write) message 1432 to AMF #1. Embedded within grant lock message 1432 is a flag indicating that AMF #1 should release the lock when finished. This optimization eliminates the need for AMF #3 to send a release lock message to AMF #1. When finished with resource X, AMF #1 sends a lock released message 1434 to AMF #3, which sends a grant lock message 1436 to AMF #4 (which is next in the queue for a write lock on resource X).

The sequence beginning with request lock (read) message 1440 shows a multiple read lock condition. With the reception of the grant lock (read) message 1442, both AMF #2 and AMF #1 have simultaneous read locks on resource Y. The write lock request 1444 sent by AMF #3 causes AMF #4 to issue release lock messages 1446 and 1448 to AMF #2 and AMF #1, respectively. This results in both lock-released message 1450 and lock-released message 1852 being sent to AMF #4. Prior to AMF #4 granting a lock to AMF #3, AMF #1 sends a request read lock message 1454 which is queued by AMF #4. AMF #3 receives the grant write lock message 1456 for resource Y which contains a flag indicating that it should release the lock when complete. AMF #3 issues a

lock released message 1458 when done with the resource Y. AMF #4 then issues a grant lock (read) message 1460 notifying AMF #1 that it has obtained a read lock on resource Y.

One benefit of this method for lock management is that the number of messages sent between the AMFs is only a function of the number of locks outstanding. It is not dependent upon the number of AMFs in the cluster configuration. Extending the first method, for example as illustrated in Figure 19, for the two-controller cluster configuration to support more than three controllers will cause the number of arbitration messages to grow with the number of AMFs in the cluster configuration. Other methods for arbitration optimization are possible, but these two methods present the preferred embodiment for the current invention.

Resource Prefetching

In the preferred embodiment for the current invention, resource arbitration is also optimized through the use of lock prefetch. An AMF can specify additional prefetch resources when arbitrating for a lock. If all or some of the prefetch resources are not locked, the Arbiter will lock them for the AMF as well. Thus, when the AMF requests the lock on these prefetched resources, (at some later time) it can quickly gain the lock (since it already had it).

20

While the invention has been described by way of example and in terms of the specific embodiments, it is to be understood that the invention is not limited to the disclosed embodiments. To the contrary, it is intended to cover various modifications and similar arrangements as would be apparent to those skilled in the art. Therefore, the scope of the appended claims should be accorded the broadest interpretation so as to encompass all such modifications and similar arrangements.

25

WHAT IS CLAIMED IS:

- 1 1. A data storage system comprising:
2 a redundancy group including a plurality of resources;
3 two or more array management functions (AMFs) sharing access to the
4 redundancy group, wherein the AMFs provide concurrent access to the redundancy group
5 for associated host systems; and
6 an interconnect medium for connecting the AMFs with the redundancy
7 group;
8 wherein when a first one of the AMFs desires to perform an operation on a
9 first resource in the redundancy group, the first AMF arbitrates with the other AMFs
10 sharing access to the redundancy group for a lock on the first resource, whereupon the
11 first AMF performs the operation on the first resource and concurrently sends replication
12 data and state information associated with the first resource to the other AMFs such that if
13 the first AMF fails while performing the operation, one of the other AMFs is able to
14 complete the operation.
- 1 2. The data storage system of claim 1, wherein the first AMF does not
2 release the lock on the first resource until a second one of the AMFs arbitrates for a lock
3 on the first resource.
- 1 3. The data storage system of claim 1, wherein if the first AMF fails,
2 the remaining AMFs arbitrate for a lock on the first resource, whereupon a second one of
3 the AMFs obtains the lock and completes the operation.
- 1 4. The data storage system of claim 1, wherein the operation
2 performed by the first AMF on the first resource includes a plurality of steps, wherein the
3 first AMF performs each step of the operation on the resource, and for each step
4 concurrently sends replication data and state information associated with the first resource
5 to the remaining AMFs, such that if the first AMF fails while performing any of the steps
6 of the operation, one of the remaining AMFs is able to complete the operation.

1 5. The data storage system of claim 1, wherein one or more host
2 systems communicate with the first AMF over one of the interconnect medium, an
3 internet connection and a PCI bus.

1 6. The system of claim 1, wherein the interconnect medium includes
2 at least one of a SCSI interface, a fiber-channel interface, a storage area network and a
3 fiber-channel storage area network.

1 7. The system of claim 1, wherein each of the AMFs executes in one
2 of a controller, a host bus adapter and a host computer.

1 8. In a data storage system, a method of dynamically sharing
2 management of a redundancy group between two or more array management functions
3 (AMFs), wherein the AMFs are able to concurrently access the redundancy group, the
4 redundancy group including a plurality of resources, the method comprising the steps of:
5 receiving a request from a host by a first one of the AMFs to perform a
6 first operation on a first one of the resources;
7 synchronizing with the other AMFs so as to acquire access to the first
8 resource; and
9 performing the first operation on the first resource.

1 9. The method of claim 8, wherein the first operation is a read
2 operation, the method further comprising the steps of:
3 receiving a request from a host by a second one of the AMFs to perform a
4 second read operation on the first resource; and
5 performing the second read operation on the first resource by the second
6 AMF concurrently with the first AMF performing the first read operation on the first
7 resource.

1 10. The method of claim 8, wherein the step of synchronizing includes
2 the step of arbitrating with the one or more other AMFs to acquire a lock on the first
3 resource so that no other AMF is able to access the first resource while the first AMF
4 maintains the lock.

1 11. The method of claim 10, further comprising the step of releasing
2 the lock on the first resource only when another AMF arbitrates for the lock.

1 12. The method of claim 8, further comprising the step of sending
2 replication data and state information associated with the first resource to the other
3 AMFs, concurrently with the step of performing the first operation, such that if the first
4 AMF fails, one of the other AMFs is able to complete the first operation.

1 13. The method of claim 12, wherein if the first AMF fails while
2 performing the first operation, the method further comprises the step of completing the
3 first operation with a second one of the AMFs.

1 14. The method of claim 13, wherein the step of completing the first
2 operation comprises the step of the second AMF synchronizing with the other AMFs so
3 as to acquire access to the first resource.

1 15. The method of claim 8, wherein the first operation includes a
2 plurality of sub-operations, wherein the step of performing the first operation includes
3 performing each sub-operation on the first resource, the method further comprising
4 sending replication data and state information associated with the first resource to the
5 other AMFs concurrently with the performance of each sub-operation, such that if the
6 first AMF fails, one of the other AMFs is able to complete the first operation.

1 16. The method of claim 8, further comprising the steps of:
2 receiving a request from a second host by a second one of the AMFs to
3 perform a second operation on a second one of the resources;
4 synchronizing with the first AMF and the other AMFs so as to acquire
5 access to the second resource; and
6 performing the second operation on the second resource.

1 17. The method of claim 16, wherein at least a portion of the second
2 operation is performed concurrently with the performance of the first operation.

1 18. A data storage network system comprising:

2 one or more redundancy groups, each redundancy group including
3 multiple resources spread over multiple disks;
4 two or more array management functions (AMFs) sharing redundancy
5 group management of the one or more redundancy groups, wherein the AMFs are able to
6 concurrently access the one or more redundancy groups; and
7 a interconnect medium for interconnecting the AMFs with the one or more
8 redundancy groups.

1 19. The system of claim 18, wherein a first one of the redundancy
2 groups includes a replacement disk, and wherein if one of the disks in the first
3 redundancy group fails, at least two of the AMFs each arbitrate for control of one or more
4 of the resources on the first redundancy group, such that each of the at least two AMFs
5 are able to concurrently reconstruct the first redundancy group using the replacement
6 disk.

1 20. The system of claim 18, wherein if an extra disk is added to a first
2 one of the redundancy groups, at least two of the AMFs each arbitrate for control of one
3 or more of the resources on the first redundancy group, such that each of the at least two
4 AMFs are able to concurrently expand the redundancy group using the extra disk.

1 21. The system of claim 18, wherein a first one of the AMFs receives a
2 write command from a host to write at least two data sets to two or more of the resources,
3 wherein the first AMF acquires a lock on the first resource to which the first data set is to
4 be written, writes the first data set to the first resource and concurrently performs a
5 replication operation wherein replication data and state information associated with the
6 first resource is sent to the other AMFs, such that if the first AMF fails while performing
7 the write operation, one of the other AMFs is able to complete the write operation.

1 22. The system of claim 21, wherein concurrently with sending
2 replication data and state information, the first AMF determines if the second data set is
3 to be written to a second resource, and if so acquires a lock on the second resource, writes
4 the second data set to the second resource and concurrently performs a second replication
5 operation wherein replication data and state information associated with the second

6 resource is sent to the other AMFs, whereupon the first AMF waits until all replication
7 operations have completed before sending status information to the host.

1 23. The system of claim 18, wherein each of the AMFs includes a
2 means for pipelining replication of incoming host data.

1 24. The system of claim 18, further comprising two or more
2 controllers, each controller implementing at least one of the two or more AMFs, wherein
3 each of the AMFs communicates with each other over one of the interconnect medium
4 and one or more PCI busses, and wherein the system further includes a means for
5 gathering messages, wherein many small messages destined for a particular one of the
6 controllers are combined and sent as one message to the particular controller.

1 25. The system of claim 18, wherein the interconnect medium includes
2 at least one of a SCSI interface, a fiber-channel interface, a storage area network and a
3 fiber-channel storage area network.

1 26. A method of reconstructing a redundancy group when one of its
2 disks fails in a data storage network system comprising two or more array management
3 functions (AMFs) interconnected with the redundancy group, wherein the redundancy
4 group includes multiple resources spread over multiple disks, wherein the redundancy
5 group also includes a replacement disk, wherein the AMFs all share management of the
6 redundancy group, and wherein the AMFs are able to concurrently access the redundancy
7 group, the method comprising the steps of:

8 arbitrating for control of a first resource by a first one of the AMFs;
9 arbitrating for control of a second resource by a second one of the AMFs;
10 and
11 concurrently reconstructing the first and second resources using the
12 replacement disk.

1 27. A method of expanding a redundancy group when an extra disk is
2 added to it in a data storage network system comprising two or more array management
3 functions (AMFs) interconnected with the redundancy group, wherein the redundancy
4 group includes multiple resources spread over multiple disks, wherein the AMFs all share

5 management of the redundancy group, and wherein the AMFs are able to concurrently
6 access the redundancy group, the method comprising the steps of:
7 arbitrating for control of a first resource by a first one of the AMFs;
8 arbitrating for control of a second resource by a second one of the AMFs;
9 and
10 concurrently expanding the first and second resources using the extra disk.

1 28. A method of pipelining replication of incoming host data in a data
2 storage network system comprising a redundancy group interconnected with two or more
3 array management functions (AMFs), wherein the redundancy group includes multiple
4 resources spread over multiple disks, wherein the AMFs all share management of the
5 redundancy group, and wherein the AMFs are able to concurrently access the redundancy
6 group, the method comprising the steps of:
7 receiving a write command by a first AMF from the host to write at least
8 two data sets to two or more of the resources;
9 acquiring a lock by the first AMF on the first resource to which the first
10 data set is to be written;
11 writing the first data set to the first resource; and concurrently
12 performing a first replication operation wherein replication data and state
13 information associated with the first resource is sent to the other AMFs, such that if the
14 first AMF fails while performing the write operation, one of the other AMFs is able to
15 complete the write operation.

1 29. The method of claim 28, further comprising the steps of:
2 concurrently with performing the first replication operation, determining if
3 the second data set is to be written to a second resource, and if so
4 acquiring a lock on the second resource;
5 writing the second data set to the second resource; and concurrently
6 performing a second replication operation wherein replication data and
7 state information associated with the second resource is sent to the other AMFs; and
8 thereafter
9 sending status information to the host after all replication operations have
10 completed.

1 30. In a data storage system, a method of dynamically sharing
2 management of a redundancy group between two or more array management functions
3 (AMFs), wherein the AMFs are able to concurrently access the redundancy group, the
4 redundancy group including a plurality of resources, the method comprising the steps of:
5 determining an arbiter AMF for a first one of the resources, wherein the
6 arbiter AMF is one of the two or more AMFs sharing management of the redundancy
7 group, and wherein the arbiter AMF is able to grant a lock for the first resource;
8 communicating a lock request from a first one of the AMFs to the arbiter
9 AMF requesting a lock on the first resource; and
10 performing an operation on the first resource by the first AMF once the
11 lock on the first resource has been granted by the arbiter AMF.

1 31. The method of claim 30, wherein a second AMF has a lock on the
2 first resource, the method further comprising the steps of:
3 issuing a release lock request to the second AMF from the arbiter AMF;
4 communicating a lock-released message from the second AMF to the
5 arbiter AMF; and thereafter
6 granting a lock on the first resource to the first AMF so that the first AMF
7 is able to perform the operation.

1 32. The method of claim 30, wherein no other AMF has a lock on the
2 first resource when the first AMF communicates the lock request, the method further
3 comprising the steps of:
4 immediately granting a lock on the first resource to the first AMF so that
5 the first AMF is able to perform the operation.

1 33. The method of claim 30, wherein the lock requested by the first
2 AMF is a write lock, and wherein once the lock is granted, no other AMF is able to obtain
3 a lock on the first resource until the first AMF releases the write lock.

1 34. The method of claim 30, wherein the lock requested by the first
2 AMF is a read lock, and wherein any of the other AMFs is able to concurrently obtain a
3 read lock on the first resource such that multiple AMFs may read the first resource
4 concurrently.

1 35. The method of claim 30, wherein the step of determining the
2 arbiter AMF includes assigning one of the two or more AMFs to be the arbiter AMF for
3 the first resource based on at least one of the number of the AMFs sharing the resource
4 and the location of the resource.

1 36. The method of claim 35, further including the step of reassigning a
2 different one of the two or more AMFs to be the arbiter AMF for the first resource if the
3 original arbiter AMF fails.

1 37. The method of claim 30, wherein the step of determining the
2 arbiter AMF includes assigning the AMF that last had a write lock on the first resource to
3 be the arbiter AMF.

1 38. The method of claim 30, wherein the arbiter AMF is able to grant a
2 lock on one or more prefetch resources, wherein the method further includes the steps of:
3 specifying a first one of the prefetch resources by the first AMF,
4 concurrently with requesting a lock on the first resource; and
5 granting a lock on the first prefetch resource to the first AMF if the first
6 prefetch is not already locked, such that the first AMF will already have the lock on the
7 first prefetch resource when it later requests the lock.

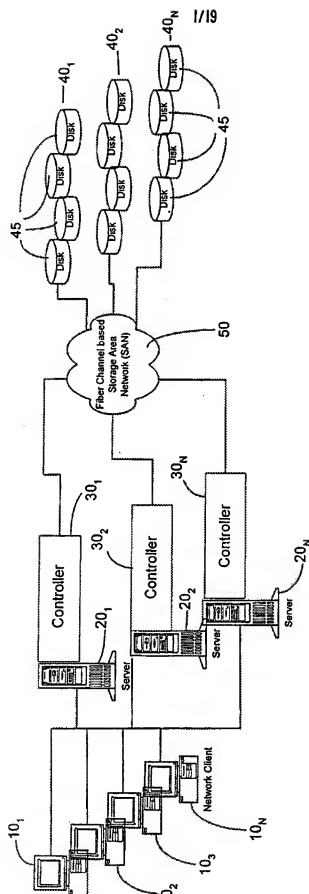
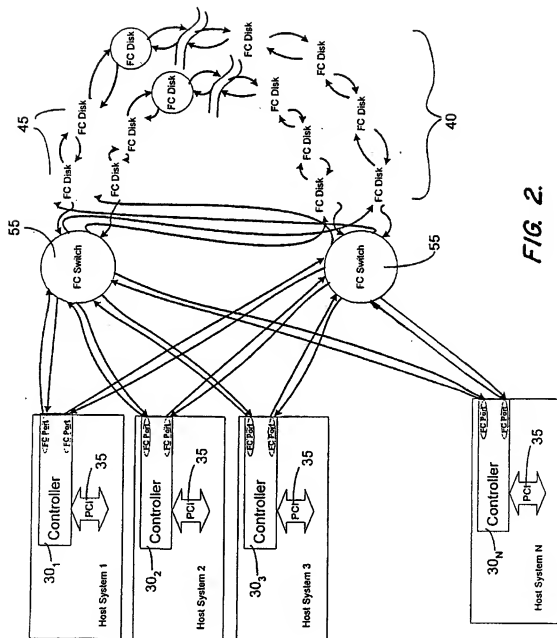


FIG. 1.

2/19



3/19

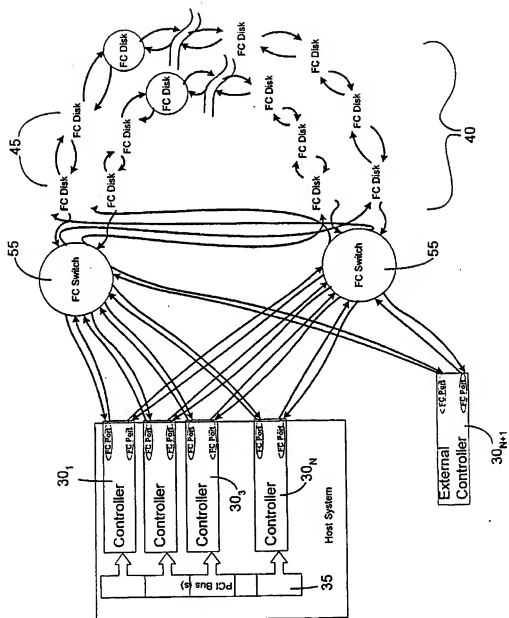
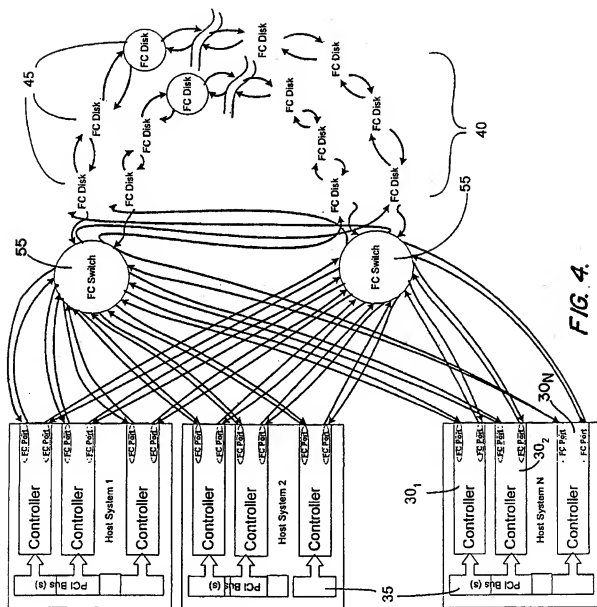


FIG. 3.

4/19



5/19

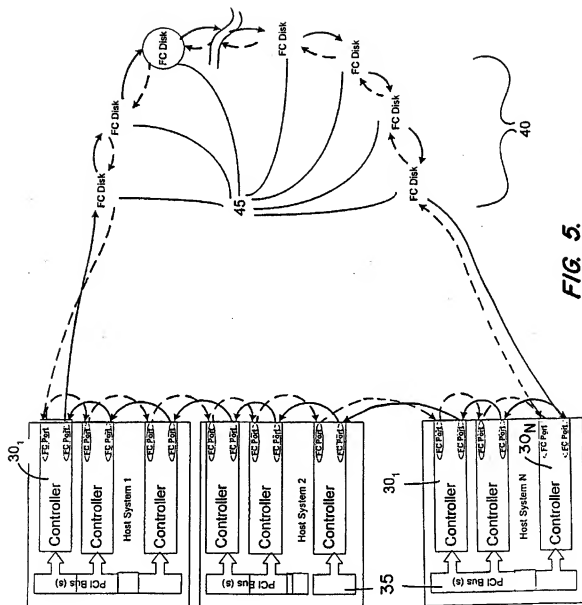
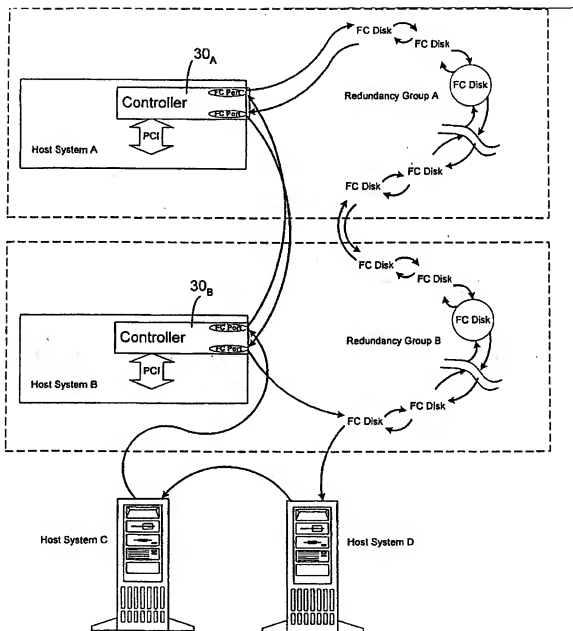


FIG. 5.

6/19

**FIG. 6.**

7/19

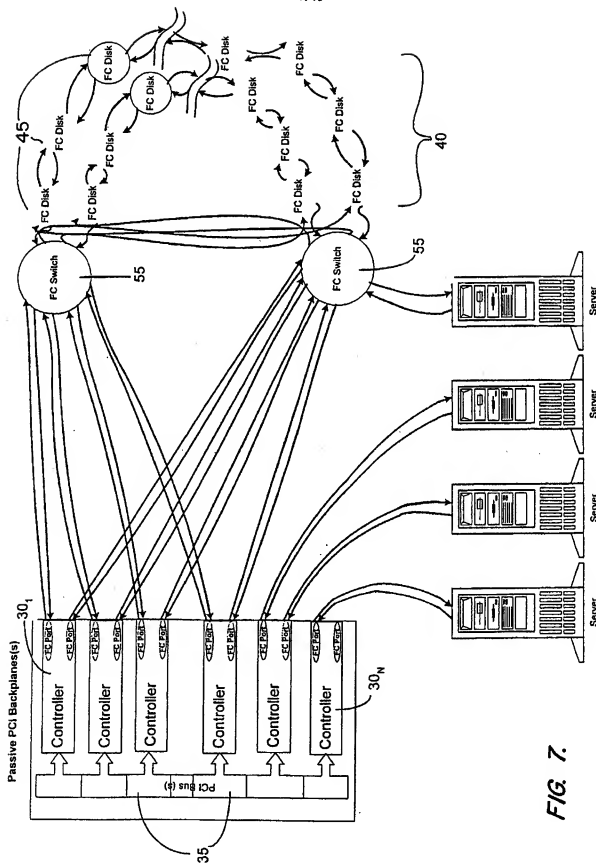


FIG. 7.

8/19

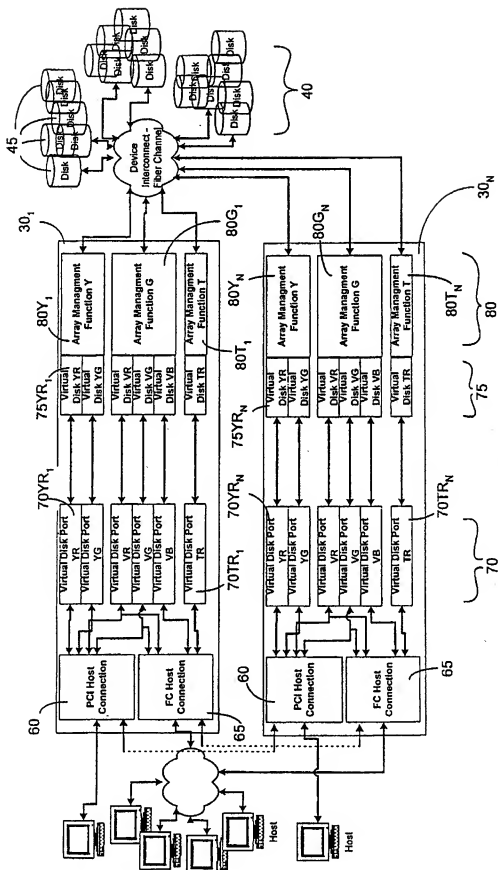


FIG. 8.

9/19

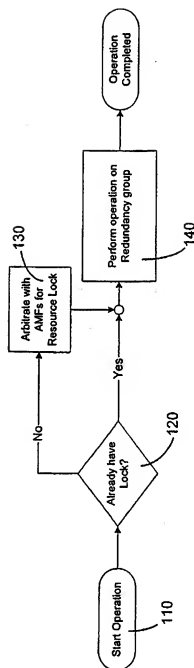


FIG. 9.

10/19

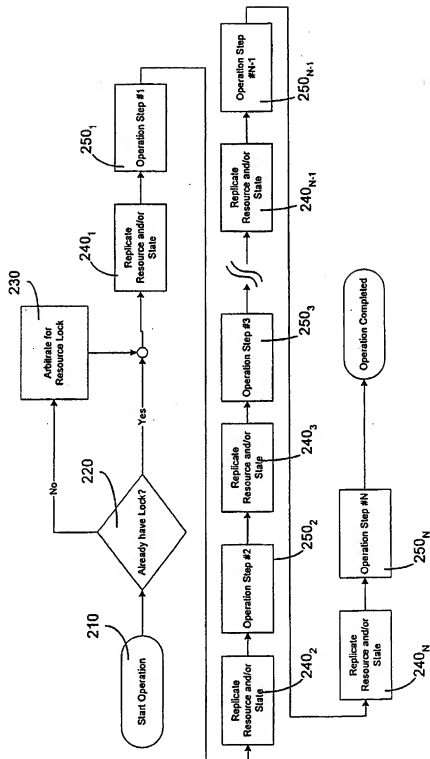
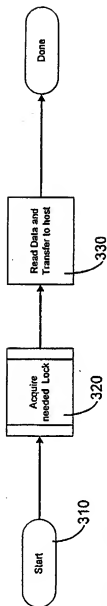
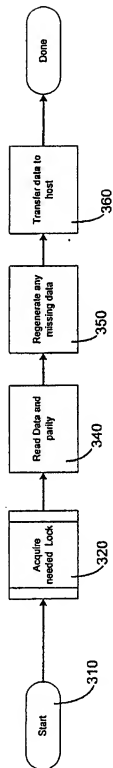


FIG. 10.

11/19

**FIG. 11a.****FIG. 11b.**

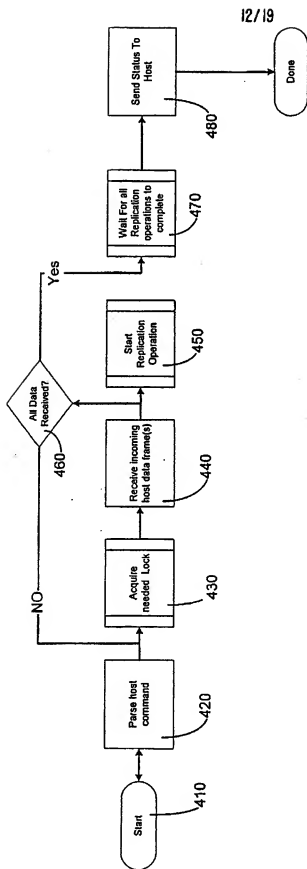
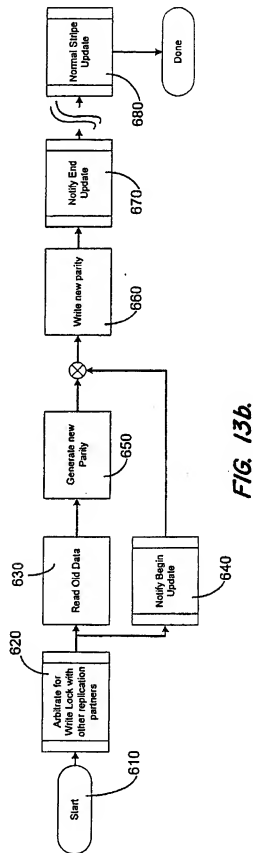
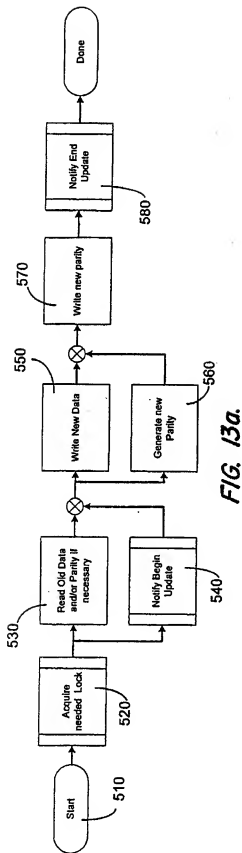
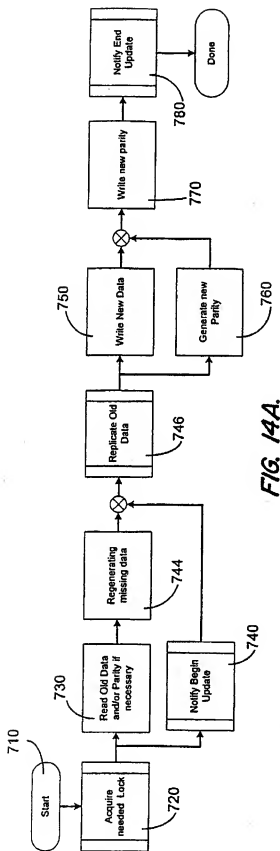
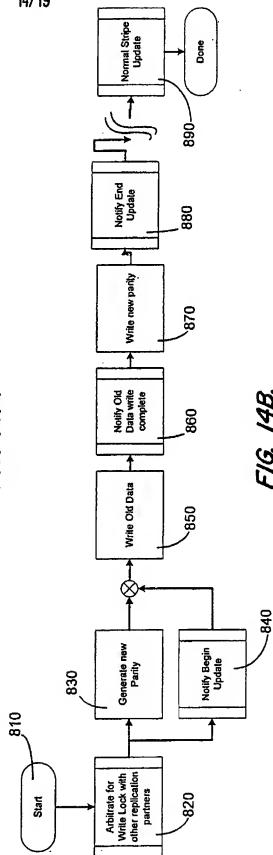


FIG. 12.





14/19



15/19

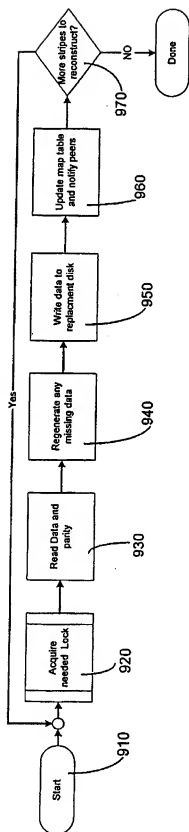


FIG. 15

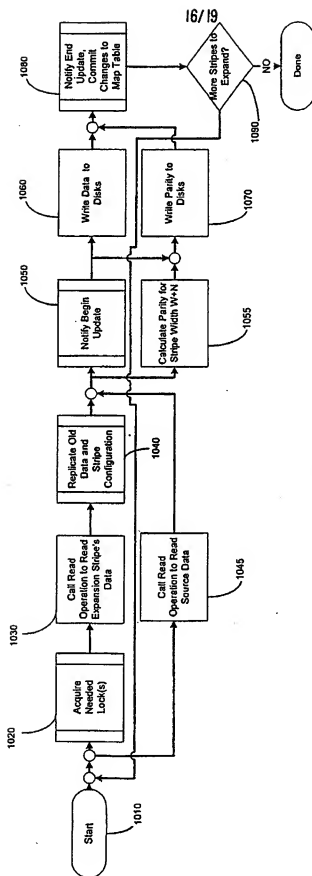


FIG. 16.

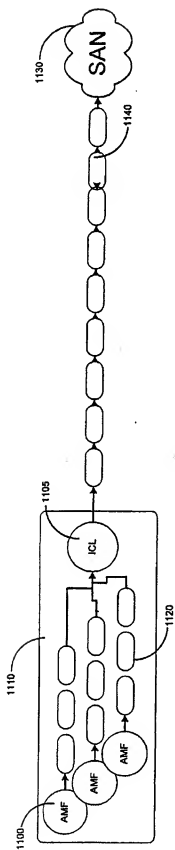


FIG. 17a.

17/19

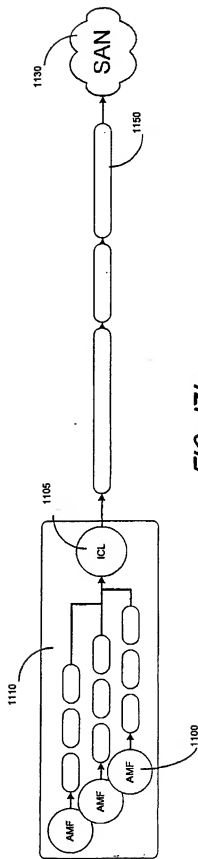


FIG. 17b.

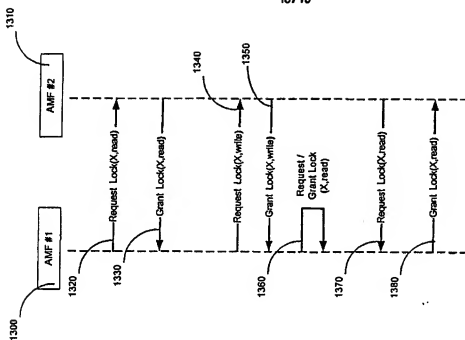


FIG. 19.

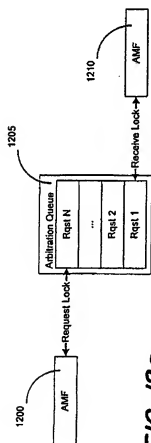


FIG. 18a.

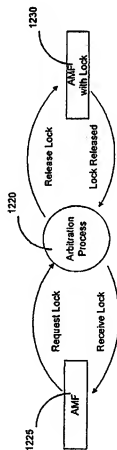


FIG. 18b.

19/19

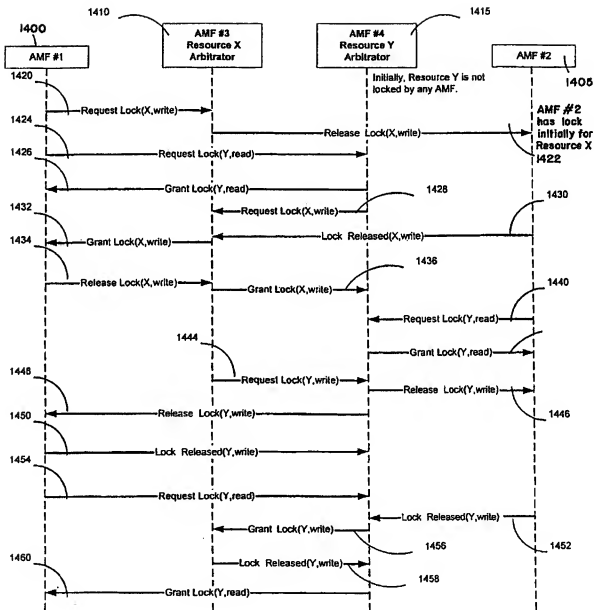


FIG. 20.

INTERNATIONAL SEARCH REPORT

International application No.
PCT/ US 00/03275

A. CLASSIFICATION OF SUBJECT MATTER

IPC(7) : G06F 11/07, 16, 20, G06F 13/14, 36; G06F 12/00
US CL : 714/9, 2, 4, 6, 7

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)

U.S. : 714/2, 4, 6, 7, 9, 27

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)

Please See Extra Sheet.

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	US 5,574,851 A (RATHUNDE) 12 November 1996, col. 1 line 1 to col. 4 line 25.	1-38
A	US 5,818,754 A (OGURA) 06 October 1998, col. 1 line 1 to col. 5 line 40.	1-38
A	US 5,694,581 A (CHENG) 02 December 1997, col. 1 line 10 to col. 4 line 63.	1-38
A	US 5,826,001 A (LUBBERS et al.) 20 October 1998, col. 1 line 5 to col. 2 line 29.	1-38

☐ Further documents are listed in the continuation of Box C. ☐ See patent family annex.

* Special categories of cited documents:	* later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention
A document defining the general state of the art which is not considered to be of particular relevance	*X* document of particular relevance; the claim of invention cannot be considered novel or obvious to be considered to involve an inventive step when the document is taken alone
E earlier document published on or after the international filing date	*Y* document of particular relevance; the claim of invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art
L document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)	*A* document member of the same patent family
O document referring to an oral disclosure, use, exhibition or other means	
P document published prior to the international filing date but later than the priority date claimed	

Date of the actual completion of the international search

05 JUNE 2000

Date of mailing of the international search report

03 JUL 2000

Name and mailing address of the ISA/US
Commissioner of Patents and Trademarks
Box PCT
Washington, D.C. 20231

Facsimile No. (703) 305-3230

Authorized officer

LY V. HUA

Telephone No. (703) 305-9684

INTERNATIONAL SEARCH REPORT

International application No.

PCT/ US 00/03275

B. FIELDS SEARCHED

Electronic data bases consulted (Name of data base and where practicable terms used):

EAST

search terms: array management functions, redundancy group, redundancy group management, arbitration, concurrent access, RAID,